



Università di Pisa

DIPARTIMENTO DI INGEGNERIA DELL'INFORMAZIONE

Laurea Triennale in Ingegneria Informatica

Implementazione del sistema di protezione contro Meltdown nel nucleo didattico

Candidato:

Riccardo Sagramoni

Matricola 565472

Relatore:

Ing. Giuseppe Lettieri

Anno Accademico 2019-2020

Indice

| | | |
|----------|---|-----------|
| 1 | Introduzione | 5 |
| 2 | Meltdown e il sistema di protezione KAISER | 7 |
| 2.1 | Background | 7 |
| 2.1.1 | Esecuzione Fuori Ordine | 7 |
| 2.1.2 | Spazi di indirizzamento | 7 |
| 2.1.3 | Attacchi Cache | 8 |
| 2.2 | Come agisce Meltdown | 8 |
| 2.2.1 | Passo 1: Leggere il segreto | 9 |
| 2.2.2 | Passo 2: Trasmettere il segreto | 9 |
| 2.2.3 | Passo 3: Ricevere il segreto | 10 |
| 2.2.4 | Conclusioni | 10 |
| 2.3 | Il sistema di protezione KAISER | 11 |
| 3 | Introduzione al nucleo didattico | 13 |
| 3.1 | Gestione dei processi | 14 |
| 3.2 | Gestione della memoria virtuale | 14 |
| 4 | Implementazione del sistema di protezione | 17 |
| 5 | Codice | 19 |
| 5.1 | sistema.cpp | 19 |
| 5.2 | sistema.s | 33 |
| 5.3 | io.s | 45 |
| 5.4 | costanti.h | 45 |
| | Bibliografia | 49 |

Capitolo 1

Introduzione

Capitolo 2

Meltdown e il sistema di protezione KAISER

La sicurezza dei sistemi informatici attuali si fonda sull'isolamento della memoria, ad esempio marcando come privilegiati gli indirizzi di memoria kernel e bloccando eventuali accessi da parte di programmi utente [14]. **Meltdown** è un tipo di attacco informatico che sfrutta un effetto collaterale dell'esecuzione fuori ordine nei processori moderni per leggere locazioni di memoria scelte in maniera arbitraria. L'attacco funziona su varie microarchitetture Intel prodotte sin dal 2010, indipendentemente dal sistema operativo in uso. Meltdown è quindi in grado di accedere arbitrariamente a qualsiasi locazione di memoria protetta (afferenti al kernel o ad altri processi) senza necessitare alcun permesso o privilegio da parte del sistema [15].

Meltdown rompe quindi tutti i meccanismi di sicurezza che si basano sull'isolamento degli spazi di indirizzamento, andando a colpire milioni di utenti. Il sistema di protezione KAISER, sviluppato originariamente per KASLR [3], ha l'importante effetto secondario di impedire l'utilizzo di Meltdown [15].

2.1 Background

2.1.1 Esecuzione Fuori Ordine

2.1.2 Spazi di indirizzamento

Per risolvere diversi problemi, in particolare l'isolamento dei processi [11], le CPU supportano l'utilizzo di spazi d'indirizzamento virtuali, in cui gli indirizzi virtuali (relativi al singolo processo) vengono tradotti in indirizzi fisici. Lo spazio d'indirizzamento di un processo (ovvero tutti i possibili indirizzi che un processo può generare) viene suddiviso in regioni dette *pagine* che possono essere mappate individualmente nella memoria fisica attraverso una tabella di traduzione multivello. Ogni processo possiede

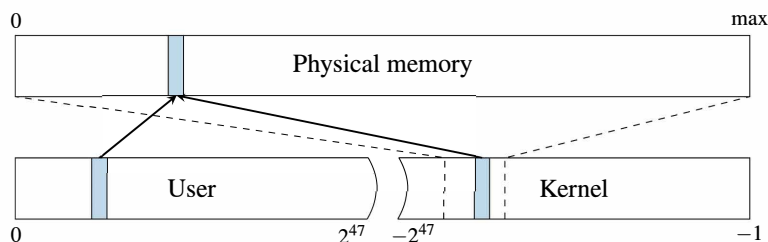


Figura 2.1: Ogni indirizzo fisico del processo è mappato sia nello spazio d'indirizzamento utente che in quello kernel all'interno della finestra di memoria fisica

una propria tabella di traduzione che traduce tutti e soli i suoi indirizzi virtuali e che definisce le proprietà di protezione delle varie zone di memoria.

Per permettere ad ogni processo di usufruire delle primitive di sistema offerte dal kernel e alle routine di sistema di accedere liberamente all'intera memoria fisica (ad esempio per modificare le tabelle di traduzione di un processo), si utilizza una traduzione, denominata *finestra di memoria fisica*, che mappa l'intera memoria fisica, compresi il kernel e lo spazio di memoria di tutti i processi, nello spazio di indirizzamento accessibile solo da livello privilegiato (figura 2.1) [13].

2.1.3 Attacchi Cache

Al fine di velocizzare gli accessi alla RAM, le CPU contengono buffer di memoria molto veloce ma di dimensioni limitate che costituiscono la cosiddetta *memoria cache*. La memoria cache maschera i tempi di latenza estremamente lunghi per l'accesso alla memoria centrale (molto lenta in confronto alla cache) conservando le locazioni di memoria che, secondo principi statistici come la *località spaziale* (se un programma accede ad un certo indirizzo, è molto probabile che in breve tempo accederà ad un indirizzo vicino) e la *località temporale* (se un programma accede ad un certo indirizzo, è molto probabile che in breve tempo vi accederà di nuovo), è più probabile vengano indirizzate dalla CPU nel breve periodo [10].

Gli attacchi a canale laterale (*side-channel attacks*) contro la cache sfruttano questa differenza di tempo di accesso introdotta dalla cache stessa. Negli attacchi Flush+Reload [16], usati da Meltdown [15], l'attaccante è in grado di determinare se una locazione di memoria è stata precedentemente caricata in cache, misurando il tempo impiegato da un'operazione di lettura.

2.2 Come agisce Meltdown

L'attacco Meltdown consiste in tre passi fondamentali [15]:

1. Leggere il contenuto di una locazione di memoria inaccessibile dall'attaccante, causando il lancio di un'eccezione di protezione
2. Accedere in maniera speculativa ad una linea di memoria cache in base al contenuto segreto della locazione protetta
3. Usare un'attacco di tipo Flush+Reload per determinare il contenuto segreto in base a quale linea di memoria è stata acceduta

2.2.1 Passo 1: Leggere il segreto

Nel primo passo di Meltdown, l'attaccante cerca di accedere ad una zona di memoria protetta, ad esempio la memoria kernel. Il tentativo di accesso ad una pagina non accessibile da livello utente fa in modo che la CPU sollevi un'eccezione di protezione, che generalmente termina il processo. Tuttavia, a causa dell'esecuzione fuori ordine, la CPU potrebbe aver già eseguito l'istruzione di accesso in maniera speculativa *prima* delle istruzioni relative all'eccezione di protezione, al fine di minimizzare i tempi di latenza (vedi paragrafo 2.1.1). In questo modo la CPU accedrebbe in maniera speculativa alla locazione di memoria desiderata prima che il processo venga terminato.

Grazie al lancio dell'eccezione, le eventuali istruzioni eseguite in maniera speculativa, che non sarebbero dovute essere eseguite in quanto relative ad una previsione di salto *errata*, non vengono *ritirate* dalla CPU e non hanno così alcun effetto sulla macroarchitettura in generale (memoria centrale e registri logici non speculativi del processore) [1].

2.2.2 Passo 2: Trasmettere il segreto

Per poter trasmettere il segreto, si utilizza un *probe array*, di dimensione pari a 256 pagine virtuali e allocato precedentemente nella memoria del processo attaccante, assicurandosi che *nessuna porzione dell'array sia presente nella cache*. La sequenza di transient instruction contiene un accesso ad un elemento del probe array, il cui offset è calcolato moltiplicando il valore del byte per la dimensione di una pagina virtuale (tipicamente e nel nostro sistema è *4KiB* [11]).

Quando la CPU gestisce l'eccezione di protezione causata dal Meltdown, le transient instruction non vengono ritirate dalla CPU, senza avere dunque effetti a livello di macroarchitettura. Sebbene quindi non sia possibile rendere direttamente disponibile il segreto dal programma utente, si hanno importanti effetti secondari a livello di microarchitettura, in particolare nella memoria cache [15].

Durante l'esecuzione speculativa, infatti, la locazione di memoria all'interno del probe array che viene acceduta dalla CPU, viene memorizzata in memoria cache e

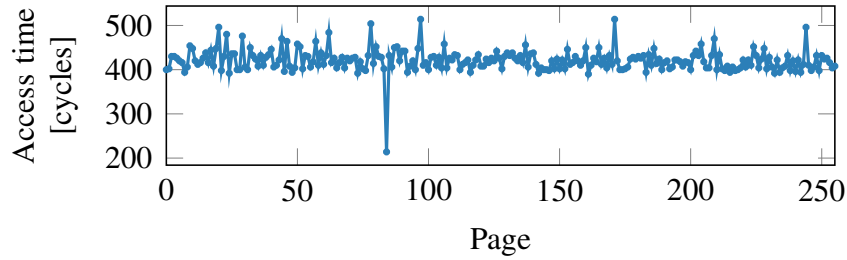


Figura 2.2: Tempo di accesso alle 256 pagine del probe array. Il grafico mostra una *cache hit* sulla pagina acceduta nel passo 2. [15]

vi rimane anche in seguito all’annullamento degli effetti delle transient instruction, rendendola vulnerabile ad un attacco side-channel.

L’utilizzo del valore segreto moltiplicato per la dimensione della pagina, ci garantisce sia una precisa *correlazione* tra il valore segreto e la locazione caricata in memoria, sia che a differenti valori della locazione di memoria saranno accedute differenti *pagine* del probe array. Ciò previene il fatto che il prefetcher hardware (per ragioni di ottimizzazione) potrebbe caricare in cache anche le locazioni di memoria adiacenti a quella acceduta, rendendo impossibile determinare quale locazione di memoria sarebbe stata indirizzata se non fosse stato utilizzato a priori questo accorgimento.

2.2.3 Passo 3: Ricevere il segreto

Dopo che la sequenza di istruzioni del passo 2 è stata eseguita, in cache è presente esattamente una linea di memoria del probe array. L’offset di questa linea è dipendente esclusivamente dal valore segreto presente nell’arbitraria locazione di memoria protetta. Grazie a ciò, l’attaccante può effettuare un’attacco Flush+Reload [16], iterando attraverso le 256 pagine del probe array e misurando il tempo di accesso per il primo elemento di ogni pagina (vedi figura 2.2). In base a quanto detto finora, la pagina con la latenza minore è l’unica presente in memoria cache e il numero della pagina è il valore segreto letto dalla memoria protetta.

2.2.4 Conclusioni

Il seguente codice mostra in Assembly x86-64 la sequenza di istruzioni alla base di Meltdown, relative ai passi 1 e 2 dell’attacco.

Listing 2.1: Meltdown frame

```
1 # rcx = indirizzo di memoria kernel
2 # rbx = probe array
3 movb (%rcx), %al          # Lettura del segreto
4 shl $12, %rax             # Traslazione dell'offset
```

```
5 movq (%rbx, %rax), %rbx # Trasmissione del segreto
```

Meltdown è quindi in grado di leggere in maniera arbitraria dati presenti in memoria *protetta*, ad esempio nello spazio di indirizzamento kernel. L'efficienza di Meltdown si basa principalmente sull'esistente *race condition* tra il lancio dell'eccezione di protezione e il passo 2 del nostro attacco (vedi paragrafo 2.2.2 a pagina 9). Per questo motivo, sono previsti alcuni accorgimenti e ottimizzazioni ulteriori non significativi per la nostra trattazione e per le quali rimandiamo a Lipp et al. [15].

Dato che l'intera memoria fisica viene mappata all'interno dello spazio di indirizzamento del kernel attraverso la cosiddetta *finestra di memoria fisica* [13], Meltdown è in grado non solo di leggere le zone di memoria relative al kernel, ma anche gli spazi di memoria di tutti gli altri processi. In base a quanto rilevato da Lipp et al. [15], Meltdown è in grado di effettuare il dump dell'intera memoria fisica fino ad una velocità di 503 KB/s.

2.3 Il sistema di protezione KAISER

KAISER, proposto da Gruss et al. [3], è una modifica del nucleo in cui il kernel non viene mappato nello spazio virtuale dei processi utente. Questa modifica era stata pensata per prevenire attacchi side-channel contro la misura di protezione KASLR [5, 4, 6], ma ha l'importante effetto secondario di prevenire Meltdown [15].

L'idea alla base di KAISER è quella di separare lo spazio di memoria kernel da quello utente, ovvero di rendere disponibile la traduzione degli indirizzi virtuali del kernel *soltanto* quando il sistema si trova in modalità privilegiata. Ciò previene Meltdown in quanto, quando il sistema sta eseguendo il programma utente della sezione 2.2.4 a fronte, l'indirizzo virtuale scelto non è presente nell'albero di traduzione del processo attaccante e quindi la CPU non è in grado di accedervi neanche in maniera speculativa. Il meccanismo proposto da Gruss et al. è quindi ritenuto la miglior soluzione a breve termine per proteggere i sistemi informatici da Meltdown [15].

KAISER introduce il concetto di **spazi d'indirizzamento shadow** per garantire l'isolamento della memoria kernel. Come mostrato in figura 2.3 nella pagina successiva, ogni processo possiede due spazi di indirizzi: uno spazio d'indirizzamento shadow, in cui è mappato solo lo spazio di memoria utente e una porzione del kernel necessaria per le interruzioni, e uno spazio d'indirizzamento in cui è mappato sia l'intero kernel che lo spazio utente, protetto con Supervisor Mode Access Prevention (SMAP) e Supervisor Mode Execution Prevention (SMEP) per compatibilità con x86 Linux [3].

Ogni qualvolta che il programma passerà da livello utente a livello sistema (ad esempio attraverso una primitiva di sistema o il lancio di un'interruzione) e viceversa, la CPU dovrà aggiornare il registro CR3 con il valore della tabella di livello 4 corrispondente

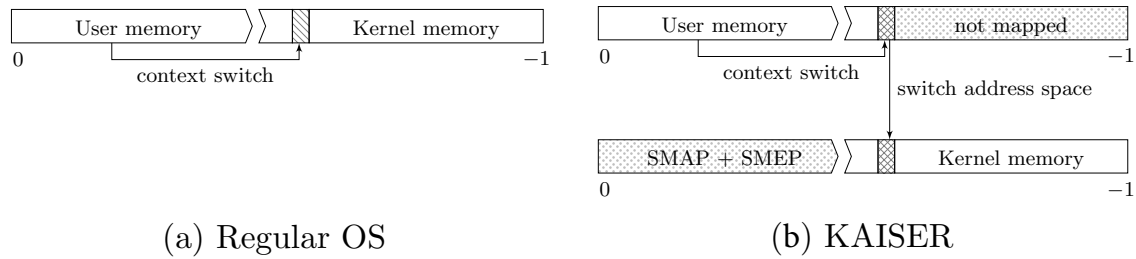


Figura 2.3: (a) Il kernel è mappato nella memoria virtuale di ogni processo. (b) TODO

al nuovo livello di privilegio (tabella shadow per il livello utente e tabella kernel per il livello sistema). Sarà dunque necessario implementare una o più funzioni (le *funzioni trampolino*), mappate in memoria shadow, dedicate all'aggiornamento corretto del registro CR3 ad ogni modifica del livello di privilegio.

Come accennato sopra, nell'architettura x86 sono necessarie alcune porzioni del kernel per il corretto funzionamento delle interruzioni e devono perciò essere mappate nello spazio di indirizzamento shadow:

- la Interrupt Descriptor Table (IDT);
- la Global Descriptor Table (GDT);
- i Task State Segment (TSS);
- la pila sistema del processo;
- la pila e il vettore delle richieste di interruzioni;
- il codice di entrata e uscita dagli interrupt handler.

Capitolo 3

Introduzione al nucleo didattico

Il nucleo didattico è un **kernel a 64 bit** perfettamente funzionante, utilizzato per finalità didattiche nel corso di Calcolatori Elettronici di Ingegneria Informatica presso l'Università di Pisa e sviluppato a partire dai concetti presentati in *Architettura dei calcolatori Vol. III* di Frosini e Lettieri[2].

Il sistema è organizzato in tre moduli distinti [8]:

- **SISTEMA**, eseguito con il processore a livello sistema, che contiene la realizzazione dei processi, inclusa la gestione della memoria;
- **IO**, eseguito con il processore a livello sistema, che contiene le routine di ingresso/uscita che permettono di utilizzare le periferiche collegate al sistema;
- **UTENTE**, eseguito con il processore a livello utente, che contiene il programma che il nucleo dovrà eseguire.

I moduli sistema e io forniscono un supporto al modulo utente, sotto forma di *primitive* che questo può invocare.

Il sistema sviluppato, per quanto funzionante, non è autosufficiente e per sviluppare i moduli necessita di un altro sistema di appoggio, nel nostro caso Linux, i cui strumenti devono essere opportunamente configurati in modo che produca eseguibili per il nostro sistema. Il nucleo così sviluppato può essere eseguito sia su una macchina reale (sconsigliato), sia su un emulatore. Nel nostro caso, useremo una versione di QEMU opportunamente modificata [9].

Il modulo sistema deve essere caricato da un *bootstrap loader* mentre il modulo io e utente devono essere caricati da una *partizione di swap*, nel nostro caso emulata da un file binario

3.1 Gestione dei processi

All'interno del nucleo didattico, i processi vengono rappresentati attraverso due strutture dati:

- Il descrittore di processo (*des_proc*), contenente il TSS del processo e i valori dei registri salvati all'ultimo cambio di contesto. Viene indirizzato dalla entrata della GDT relativa al processo.
- Il *proc_elem*, contenente l'id e la priorità del processo e usato nelle code di processi.

Il sistema prevede una politica di schedulazione a priorità fissa, in cui passa in esecuzione il processo pronto con priorità *maggiore*. A parità di priorità, viene adottata una politica FIFO (*First Input, First Output*).

3.2 Gestione della memoria virtuale

Il nostro sistema implementa la paginazione su domanda [12], con zone di memoria condivise tra tutti i processi e zone private ai processi [7].

La memoria virtuale di ogni processo è diviso nelle seguenti sezioni (vedi anche 3.1:

- Il sottospazio canonico superiore, con indirizzi da 0x0000000000000000 a 0x00007FFFFFFFFFFFFF, accessibile solo da livello sistema. A sua volta suddiviso in:
 - **sistema/condivisa**: contiene la finestra di memoria fisica
 - **sistema/privata**: contiene la pila sistema del processo
 - **IO/condivisa**: contiene il modulo I/O
- Il sottospazio canonico inferiore, con indirizzi da 0xFFFFF80000000000 a 0xFFFFFFFFFFFFFFFF, accessibile solo da livello sistema. A sua volta suddiviso in:
 - **utente/condivisa**: contiene il modulo utente, ovvero le sezioni `.text` e `.data` del programma utente
 - **utente/privata**: contiene la pila utente del processo

Si noti come alcune parti della memoria fisica siano accessibili esclusivamente tramite la finestra: il primo MiB di memoria (riservato per ragioni storiche), lo heap di sistema, il modulo sistema, i descrittori di frame e i descrittori di pagine virtuali.

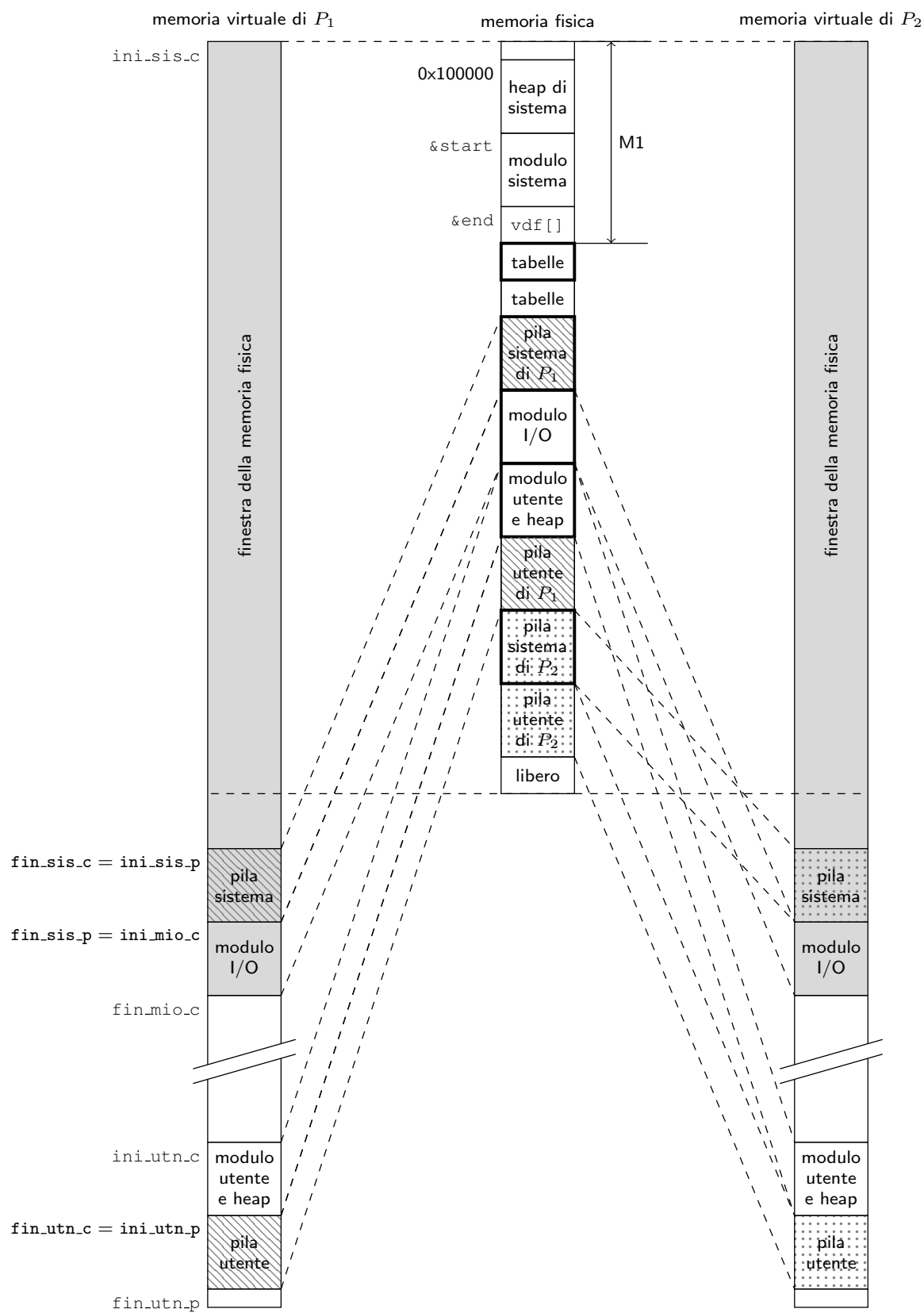


Figura 3.1: Esempio di memoria virtuale con due processi (non in scala) [7]

Usando la paginazione su domanda, la memoria fisica funziona da cache dello swap. Normalmente le sole sezioni *residenti*, ovvero non selezionabili come vittime per uno swap, sarebbero la *sistema/condivisa* e la *sistema/privata*, in quanto contenenti la finestra di memoria e la pila usata dal modulo sistema. Per semplicità di implementazione, nel nostro sistema sono state rese residenti tutte le sezioni condivise [7].

Capitolo 4

Implementazione del sistema di protezione

Capitolo 5

Codice

5.1 sistema.cpp

Listing 5.1: Descrittori di processo e di TSS

```
19 // TASK STATE SEGMENT, RICHIESTO DALL'HARDWARE.
20 // PER PROTEGGERE IL SISTEMA DA MELTDOWN, E' UNICO PER TUTTI
21 // I PROCESSI ED E' MAPPATO NELLA MEMORIA SHADOW DEL KERNEL
22 struct __attribute__((packed)) task_state_segment {
23     // PARTE RICHIESTA DALL'HARDWARE
24     natl riservato1;
25     vaddr punt_nucleo;
26     // DUE QUAD A DISPOSIZIONE (PUNTATORI ALLE PILE RING 1 E
27     ↪ 2)
28     natq disp1[2];
29     natq riservato2;
30     //ENTRY DELLA IST, NON USATA
31     natq disp2[7];
32     natq riservato3;
33     natw riservato4;
34     natw iomap_base; // SI VEDA INIT_TSS()
35 };
36 extern "C" task_state_segment tss;
37
38 // DESCRITTORE SOFTWARE DI PROCESSO
39 struct des_proc {
40     vaddr punt_nucleo; // PUNTATORE ALLA PILA SISTEMA DEL
41     ↪ PROCESSO
42     faddr cr3;
```

```

42     natq contesto[N_REG];    // ARRAY PER SALVARE IL CONTESTO
43     natl cpl;
44 };
45
46 // ARRAY DEI DESCRITTORI DI PROCESSO
47 des_proc* descrittori_processi[NUM_TSS];
48 // POSIZIONE IMMEDIATAMENTE SUCCESSIVA ALL'ULTIMO ID
    ↪ ASSEGNATO
49 natl next_id = 0;
50
51 // DATO UN ID, RESTITUISCE IL PUNTATORE AL CORRISPONDENTE
    ↪ DES_PROC
52 extern "C" des_proc* des_p(natl id)
53 {
54     return descrittori_processi[id];
55 }
56
57 // CERCA UN'ENTRATA LIBERA NELL'ARRAY DESCRITTORI_PROCESSI DA
    ↪ ASSEGNARE
58 // AL PROCESSO IL CUI DESCRITTORE E' PASSATO COME PARAMETRO.
59 // RESTITUISCE L'OFFSET DELL'ENTRATA ASSEGNATA (L'ID DEL
    ↪ NUOVO PROCESSO)
60 // OPPURE NUM_TSS SE NON CI SONO PIU' ID LIBERI
61 natl seleziona_id_processo (des_proc* d)
62 {
63     natl i = next_id;
64
65     do {
66         if (descrittori_processi[i] == 0) {
67             // TROVATO DESCRITTORE DI PROCESSO LIBERO
68             descrittori_processi[i] = d;
69             next_id = (i + 1) % NUM_TSS;
70             return i;
71         }
72
73         i = (i + 1) % NUM_TSS;
74     } while (i != next_id);
75
76     // NON CI SONO PIU' ID LIBERI
77     return NUM_TSS;
78 }

```

```

79
80 // RENDE LIBERO L'ID PASSATO COME PARAMETRO, AZZERANDO LA
    ↳ RELATIVA
81 // ENTRATA NELL'ARRAY DESCRITTORI PROCESSI
82 void rilascia_id_processo (natl id)
83 {
84     descrittori_processi[id] = 0;
85 }

```

Listing 5.2: Creazione della finestra di memoria shadow

```

743 // CONFINI SEZIONI TRAMPOLINO
744 extern "C" natq start_trampoline_text;
745 extern "C" natq end_trampoline_text;
746 extern "C" natq start_trampoline_data;
747 extern "C" natq end_trampoline_data;
748 extern "C" natq start_trampoline_bss;
749 extern "C" natq end_trampoline_bss;
750
751 // TAB_ENTRY CHE DESCRIVE LA FINESTRA DI MEMORIA SHADOW
752 // DA SOSTITUIRE CON LA FINESTRA DI MEMORIA NORMALE
753 // QUANDO IL PROGRAMMA PASSA A LIVELLO UTENTE
754 tab_entry des_finestra_shadow = 0;
755
756 // TAB_ENTRY CHE DESCRIVE LA FINESTRA DI MEMORIA KERNEL
    ↳ COMPLETA.
757 extern "C" tab_entry des_finestra_kernel;
758
759 void mappa_pagina_shadow (vaddr ind_virt, faddr tab3, bool
    ↳ isText);
760
761 // RESTITUISCE UN RIFERIMENTO AL DESCRITTORE DI LIVELLO LIV
762 // DA CUI PASSA LA TRADUZIONE DELL'INDIRIZZO IND_VIRT NELLA
    ↳ MEMORIA SHADOW
763 tab_entry& get_des_shadow (int livello, vaddr ind_virt, faddr
    ↳ tab3)
764 {
765     faddr tab = tab3;
766
767     for (int i = 3; i > livello; i--) {
768         // PRELEVA IL DESCRITTORE DELLA TABELLA DI LIVELLO I-
            ↳ ESIMO

```

```

769         tab_entry entry = get_entry(tab, i_tab(ind_virt, i));
770
771         if (!extr_P(entry))
772             panic("P=0_non_ammesso");
773
774         tab = extr_IND_FISICO(entry);
775     }
776
777     return get_entry(tab, i_tab(ind_virt, livello));
778 }
779
780 // SE ASSENTE, CREA UNA NUOVA TABELLA DI TRADUZIONE DI
781     ↪ LIVELLO 1-3
782 // O MAPPA LA TRADUZIONE DI UNA PAGINA VIRTUALE AL LIVELLO 1
783 // NELLO SPAZIO DI TRADUZIONE SHADOW
784 void crea_shadow (vaddr ind_virt, int liv, faddr tab3, bool
785     ↪ isText)
786 {
787     if (liv < 0 || liv > 3) {
788         panic("crea_shadow(...):_valore_liv_non_valido");
789     }
790
791     tab_entry& dt = get_des_shadow(liv + 1, ind_virt, tab3);
792     bool bit_P = extr_P(dt);
793
794     if (!bit_P && liv != 0) {
795         des_frame* df = alloca_frame_libero();
796
797         if (df == 0) {
798             flog(LOG_ERR, "Impossibile_allocare_copia_shadow_
799                 ↪ di_una_pagina");
800             panic("errore");
801         }
802
803         // INIZIALIZZA DESCRITTORE DI PAGINA FISICA
804         df->livello = liv;
805         df->residente = true;
806         df->processo = esecuzione->id;
807         df->ind_virtuale = ind_virt;
808
809         // INIZIALIZZA NUOVO DESCRITTORE DI TABELLA O PAGINA

```

```

807         ↪ VIRTUALE
808         faddr new_entry = indirizzo_frame(df);
809         memset(reinterpret_cast<void*>(new_entry), 0,
810             ↪ DIM_PAGINA);
811
812         // COLLEGA IL NUOVO DESCRITTORE AL PRECEDENTE
813         set_IND_FISICO(dt, new_entry);
814         set_P(dt, true);
815         dt |= BIT_RW;
816     }
817     else if (!bit_P) {
818         // MAPPA UNA PAGINA VIRTUALE
819         set_IND_FISICO(dt, ind_virt);
820         set_P(dt, true);
821         if (!isText) dt = dt | BIT_RW;
822     }
823 }
824
825 void mappa_pagina_shadow (vaddr ind_virt, faddr tab3, bool
826     ↪ isText)
827 {
828     // CREA LE TABELLE DI TRADUZIONE DAL LIVELLO 2 IN POI
829     for (int i = 2; i >= 0; i--) {
830         crea_shadow(ind_virt, i, tab3, isText);
831     }
832 }
833
834 // CREA LA TABELLA DI LIVELLO 3 "SHADOW", OVVERO LA FM
835     ↪ COSTITUITA
836 // DALLE PAGINE KERNEL STRETTAMENTE NECESSARIO AI PROCESSI
837     ↪ UTENTI
838 // E LE FUNZIONI TRAMPOLINO PER PASSARE NELLA FM COMPLETA (ED
839     ↪ USCIRVI)
840 faddr crea_tab3_shadow ()
841 {
842     des_frame* df = alloca_frame_libero();
843     if (df == 0) {
844         flog(LOG_ERR, "Impossibile_allocare_copia_shadow_
845             ↪ della_finestra_di_memoria");
846         panic("errore");
847     }
848 }

```

```

841
842 // INIZIALIZZA IL DESCRITTORE DI FRAME
843 df->livello = 3;
844 df->residente = true;
845 df->processo = esecuzione->id;
846 df->ind_virtuale = 0;
847
848 // INIZIALIZZA LA TABELLA
849 faddr tab3 = indirizzo_frame(df);
850 memset(reinterpret_cast<void*>(tab3), 0, DIM_PAGINA);
851
852 return tab3;
853 }
854
855 // MAPPA NELLA MEMORIA VIRTUALE SHADOW
856 // LA SEZIONE TEXT, DATA E BSS DELLA FINESTRA DI MEMORIA
857 // ↳ TRAMPOLINO
858 // CHE SONO ALLINEATE A 4KiB
859 void mappa_modulo_sistema_trampoline_in_shadow (faddr tabFM)
860 {
861     natq dim_trampoline_text = (natq)&end_trampoline_text - (
862         ↳ natq)&start_trampoline_text;
863     natq num_pag_text = (dim_trampoline_text / DIM_PAGINA) +
864         ↳ 1;
865
866     for (natq i = 0; i < num_pag_text; i++) {
867         mappa_pagina_shadow((natq)&start_trampoline_text + i*
868             ↳ DIM_PAGINA, tabFM, true);
869     }
870
871     natq dim_trampoline_data = (natq)&end_trampoline_data - (
872         ↳ natq)&start_trampoline_data;
873     natq num_pag_data = (dim_trampoline_data / DIM_PAGINA) +
874         ↳ 1;
875
876     for (natq i = 0; i < num_pag_data; i++) {
877         mappa_pagina_shadow((natq)&start_trampoline_data + i*
878             ↳ DIM_PAGINA, tabFM, false);
879     }
880
881     natq dim_trampoline_bss = (natq)&end_trampoline_bss - (

```



```

    ↪ natq)&start_trampoline_bss;
875     natq num_pag_bss = (dim_trampoline_bss / DIM_PAGINA) + 1;
876
877     for (natq i = 0; i < num_pag_bss; i++) {
878         mappa_pagina_shadow((natq)&start_trampoline_bss + i*
            ↪ DIM_PAGINA, tabFM, false);
879     }
880 }
881
882 void crea_finestra_FM_shadow (faddr tab4)
883 {
884     faddr tab3_shadow = crea_tab3_shadow();
885
886     set_IND_FISICO(des_finestra_shadow, tab3_shadow);
887     set_P(des_finestra_shadow, true);
888     des_finestra_shadow |= BIT_RW;
889
890     // MAPPA LA SEZIONE TEXT, DATA E BSS DELLA FINESTRA DI
        ↪ MEMORIA TRAMPOLINO
891     mappa_modulo_sistema_trampoline_in_shadow(tab3_shadow);
892 }
893
894 // MAPPA LA MEMORIA FISICA IN MEMORIA VIRTUALE, INCLUSA L'
        ↪ AREA PCI
895 // (COPIAMO LA FINESTRA GIA' CREATA DAL BOOT LOADER)
896 bool crea_finestra_FM(faddr tab4)
897 {
898     faddr boot_dir = readCR3();
899     copy_des(boot_dir, tab4, I_SIS_C, N_SIS_C);
900
901     crea_finestra_FM_shadow(tab4);
902
903     // SALVA LA TAB_ENTRY RELATIVA ALLA FINESTRA DI MEMORIA
        ↪ IN
904     // MODALITA' KERNEL, NELLA MEMORIA TRAMPOLINO
905     des_finestra_kernel = *(reinterpret_cast<tab_entry*>(tab4
        ↪ ));
906
907     return true;
908 }

```

Listing 5.3: Creazione di un processo

```

1009 proc_elem* crea_processo(void f(int), int a, int prio, char
    ↪ liv, bool IF)
1010 {
1011     proc_elem* p; // PROC_ELEM PER IL NUOVO PROCESSO
1012     natl identifier; // IDENTIFICATORE DEL PROCESSO
1013     des_proc* pdes_proc; // DESCRITTORE DI PROCESSO
1014     des_frame* dpf_tab4; // TAB4 DEL PROCESSO
1015     faddr pila_sistema;
1016
1017     // ( ALLOCAZIONE (E AZZERAMENTO PREVENTIVO) DI UN
    ↪ DES_PROC
1018     pdes_proc = static_cast<des_proc*>(alloca(sizeof(des_proc
    ↪ ))));
1019     if (pdes_proc == 0) goto errore1;
1020     memset(pdes_proc, 0, sizeof(des_proc));
1021     // )
1022
1023     // ( SELEZIONE DI UN IDENTIFICATORE
1024     identifier = seleziona_id_processo(pdes_proc);
1025     if (identifier == NUM_TSS) goto errore2;
1026     // )
1027
1028     // ( ALLOCAZIONE E INIZIALIZZAZIONE DI UN PROC_ELEM
1029     p = static_cast<proc_elem*>(alloca(sizeof(proc_elem)));
1030     if (p == 0) goto errore3;
1031     p->id = identifier;
1032     p->precedenza = prio;
1033     p->puntatore = 0;
1034     // )
1035
1036     // ( CREAZIONE DELLA TAB4 DEL PROCESSO
1037     dpf_tab4 = alloca_frame(p->id, 4, 0);
1038     if (dpf_tab4 == 0) goto errore4;
1039     dpf_tab4->livello = 4;
1040     dpf_tab4->residente = true;
1041     dpf_tab4->processo = identifier;
1042     pdes_proc->cr3 = indirizzo_frame(dpf_tab4);
1043     crea_tab4(pdes_proc->cr3);
1044     mappa_pagina_shadow(pdes_proc->cr3, extr_IND_FISICO(
    ↪ des_finestra_shadow), false);

```

```

1045 // )
1046
1047 // ( CREAZIONE DELLA PILA SISTEMA .
1048 if (!crea_pila(p->id, fin_sis_p, DIM_SYS_STACK,
    ↪ LIV_SISTEMA))
1049     goto errore5;
1050 pila_sistema = carica_pila_sistema(p->id, fin_sis_p,
    ↪ DIM_SYS_STACK);
1051 if (pila_sistema == 0)
1052     goto errore6;
1053 // )
1054
1055 if (liv == LIV_UTENTE) {
1056     // ( INIZIALIZZIAMO LA PILA SISTEMA .
1057     natq* pl = reinterpret_cast<natq*>(pila_sistema);
1058
1059     pl[-5] = reinterpret_cast<natq>(f); // RIP (CODICE
    ↪ UTENTE)
1060     pl[-4] = SEL_CODICE_UTENTE; // CS (CODICE UTENTE)
1061     pl[-3] = IF ? BIT_IF : 0; // RFLAGS
1062     pl[-2] = fin_utn_p - sizeof(natq); // RSP
1063     pl[-1] = SEL_DATI_UTENTE; // SS (PILA UTENTE)
1064     // ESEGUENDO UNA IRET DA QUESTA SITUAZIONE, IL
    ↪ PROCESSO
1065     // PASSERA' AD ESEGUIRE LA PRIMA ISTRUZIONE DELLA
    ↪ FUNZIONE F,
1066     // USANDO COME PILA LA PILA UTENTE (AL SUO
    ↪ INDIRIZZO VIRTUALE)
1067 // )
1068
1069 // ( CREAZIONE DELLA PILA UTENTE
1070 if (!crea_pila(p->id, fin_utn_p, DIM_USR_STACK,
    ↪ LIV_UTENTE)) {
1071     flog(LOG_WARN, "creazione_pila_utente_fallita");
1072     goto errore6;
1073 }
1074 // )
1075
1076 // ( INFINE, INIZIALIZZIAMO IL DESCRITTORE DI
    ↪ PROCESSO
1077 // INDIRIZZO DEL BOTTOM DELLA PILA SISTEMA, CHE

```

```

1078 // DAL MECCANISMO DELLE INTERRUZIONI
1079 pdes_proc->punt_nucleo = fin_sis_p;
1080
1081 // INIZIALMENTE, IL PROCESSO SI TROVA A LIVELLO
1082 // SE AVESSE ESEGUITO UNA ISTRUZIONE INT, CON LA
1083 // CHE CONTIENE LE 5 PAROLE LUNGHE PREPARATE
1084 pdes_proc->contesto[I_RSP] = fin_sis_p - 5 * sizeof(
1085 // IL REGISTRO RDI DEVE CONTENERE IL PARAMETRO DA
1086 // PASSARE
1087 // ALLA FUNZIONE F
1088 pdes_proc->contesto[I_RDI] = a;
1089 //PDES_PROC->CONTESTO[I_FPU_CR] = 0x037F;
1090 //PDES_PROC->CONTESTO[I_FPU_TR] = 0xFFFF;
1091 pdes_proc->cpl = LIV_UTENTE;
1092
1093 // TUTTI GLI ALTRI CAMPI VALGONO 0
1094 // )
1095 } else {
1096 // ( INIZIALIZZAZIONE DELLA PILA SISTEMA
1097 natq* pl = reinterpret_cast<natq*>(pila_sistema);
1098 pl[-6] = reinterpret_cast<natq>(f); // RIP (CODICE
1099 // SISTEMA)
1100 pl[-5] = SEL_CODICE_SISTEMA; // CS (CODICE SISTEMA)
1101 pl[-4] = IF ? BIT_IF : 0; // RFLAGS
1102 pl[-3] = fin_sis_p - sizeof(natq); // RSP
1103 pl[-2] = 0; // SS
1104 pl[-1] = 0; // IND. RIT.
1105 // (NON SIGNIFICATIVO)
1106 // I PROCESSI ESTERNI LAVORANO ESCLUSIVAMENTE A
1107 // LIVELLO
1108 // SISTEMA. PER QUESTO MOTIVO, PREPARIAMO UNA SOLA
1109 // PILA (LA
1110 // PILA SISTEMA)
1111 // )

```

```

1110      // ( INIZIALIZZIAMO IL DESCRITTORE DI PROCESSO
1111      pdes_proc->contesto[I_RSP] = fin_sis_p - 6 * sizeof(
           ↪ natq);
1112      pdes_proc->contesto[I_RDI] = a;
1113
1114      //PDES_PROC->CONTESTO[I_FPU_CR] = 0x037F;
1115      //PDES_PROC->CONTESTO[I_FPU_TR] = 0xFFFF;
1116      pdes_proc->cpl = LIV_SISTEMA;
1117
1118      //  TUTTI GLI ALTRI CAMPI VALGONO 0
1119      // )
1120  }
1121
1122  return p;
1123
1124 errore6:    rilascia_tutto(indirizzo_frame(dpf_tab4), I_SIS_P
           ↪ , N_SIS_P);
1125 errore5:    rilascia_frame(dpf_tab4);
1126 errore4:    dealloca(p);
1127 errore3:    rilascia_id_processo(identifier);
1128 errore2:    dealloca(pdes_proc);
1129 errore1:    return 0;
1130 }

```

Listing 5.4: Distruzione di un processo

```

1191 void distruggi_processo(proc_elem* p)
1192 {
1193     des_proc* pdes_proc = des_p(p->id);
1194
1195     faddr tab4 = pdes_proc->cr3;
1196     riassegna_tutto(p->id, tab4, I_MIO_C, N_MIO_C);
1197     riassegna_tutto(p->id, tab4, I_UTN_C, N_UTN_C);
1198     rilascia_tutto(tab4, I_UTN_P, N_UTN_P);
1199     ultimo_terminato = tab4;
1200     if (p != esecuzione) {
1201         distruggi_pila_precedente();
1202     }
1203     rilascia_id_processo(p->id);
1204     dealloca(pdes_proc);
1205 }

```

Listing 5.5: Inizializzazione del sistema

```

1838 // AZZERRA TUTTE LE ENTRATE DELL'ARRAY DESCRITTORI_PROCESSI,
1839 // RENDENDO LIBERI TUTTI GLI ID DISPONIBILI
1840 void init_descrittori_processi ()
1841 {
1842     memset(descrittori_processi, 0, sizeof(des_proc*) *
           ↪ NUM_TSS);
1843 }
1844
1845 void init_tss ()
1846 {
1847     // IL CAMPO IOMAP_BASE CONTIENE L'OFFSET (NEL TSS) DELL
           ↪ 'INIZIO
1848     // DELLA "I/O BITMAP". QUESTA BITMAP CONTIENE UN BIT
           ↪ PER OGNI
1849     // POSSIBILE INDIRIZZO DI I/O. LE ISTRUZIONI IN E OUT
           ↪ ESEGUITE
1850     // DA LIVELLO UTENTE VERRANNO PERMESSE SE IL BIT
           ↪ CORRISPONDENTE
1851     // ALL'INDIRIZZO DI I/O A CUI SI RIFERISCONO VALE 1.
1852     // PER DISATTIVARE QUESTO MECCANISMO DOBBIAMO
           ↪ INIZIALIZZARE
1853     // IL CAMPO IOMAP_BASE CON UN OFFSET MAGGIORE O UGUALE
1854     // DELLA DIMENSIONE DEL SEGMENTO TSS (COME SCRITTA NEL
1855     // DESCRITTORE DI SEGMENTO TSS NELLA GDT, VEDERE '
           ↪ SET_ENTRY_TSS'
1856     // IN SISTEMA.S)
1857     tss.iomap_base = DIM_DESP;
1858
1859     // GLI ALTRI CAMPI DEL TSS SONO NULLI
1860 }
1861
1862 extern "C" void salta_a_main();
1863 extern "C" void cmain()
1864 {
1865     natl mid;
1866
1867     // (* ANCHE SE IL PRIMO PROCESSO NON E' COMPLETAMENTE
           ↪ INIZIALIZZATO ,
1868     // GLI Diamo UN IDENTIFICATORE, IN MODO CHE COMPAIA
           ↪ NEI LOG

```

```

1869     init.id = 0xFFFFFFFF;
1870     init.precedenza = MAX_PRIORITY;
1871     esecuzione = &init;
1872     // *)
1873
1874     flog(LOG_INFO, "Nucleo_di_Calcolatori_Elettronici,_v5
        ↪ .12.6_con_patch_contro_Meltdown");
1875     init_tss();
1876     flog(LOG_INFO, "tss_inizializzato");
1877     init_gdt();
1878     flog(LOG_INFO, "gdt_inizializzata");
1879
1880     // (* ASSEGNA ALLO HEAP DI SISTEMA HEAP_SIZE BYTE NEL
        ↪ SECONDO MiB
1881     heap_init((addr)HEAP_START, HEAP_SIZE);
1882     flog(LOG_INFO, "Heap_di_sistema:_%x_B_@%x", HEAP_SIZE,
        ↪ HEAP_START);
1883     // *)
1884
1885     // ( IL RESTO DELLA MEMORIA E' PER I FRAME (PARTE M2)
1886     init_des_frame();
1887     flog(LOG_INFO, "Pagine_fisiche:_%d", N_DF);
1888     // )
1889
1890     flog(LOG_INFO, "sis/cond_[%p,_%p)", ini_sis_c, fin_sis_c)
        ↪ ;
1891     flog(LOG_INFO, "sis/priv_[%p,_%p)", ini_sis_p, fin_sis_p)
        ↪ ;
1892     flog(LOG_INFO, "io_/cond_[%p,_%p)", ini_mio_c, fin_mio_c)
        ↪ ;
1893     flog(LOG_INFO, "usr/cond_[%p,_%p)", ini_utn_c, fin_utn_c)
        ↪ ;
1894     flog(LOG_INFO, "usr/priv_[%p,_%p)", ini_utn_p, fin_utn_p)
        ↪ ;
1895
1896     faddr inittab4 = crea_tab4();
1897
1898     if(!crea_finestra_FM(inittab4))
1899         goto error;
1900     loadCR3(inittab4);
1901     flog(LOG_INFO, "Caricato_CR3");

```

```

1902
1903     apic_init(); // IN LIBCE
1904     apic_reset(); // IN LIBCE
1905     apic_fill();
1906     flog(LOG_INFO, "APIC_inizializzato");
1907
1908     // ( INIZIALIZZAZIONE DELLO SWAP, CHE COMPRENDE LA
1909         ↳ LETTURA
1910     //     DEGLI ENTRY POINT DI START_IO E START_UTENTE
1911     if (!swap_init())
1912         goto error;
1912     flog(LOG_INFO, "sb:_blocks_=_%d", swap_dev.sb.blocks);
1913     flog(LOG_INFO, "sb:_user_=_%p/%p",
1914         swap_dev.sb.user_entry,
1915         swap_dev.sb.user_end);
1916     flog(LOG_INFO, "sb:_io_=_%p/%p",
1917         swap_dev.sb.io_entry,
1918         swap_dev.sb.io_end);
1919     // )
1920
1921     // ( INIZIALIZZA L'ARRAY CHE CONTIENE I PUNTATORI AI
1922         ↳ DESCRITTORI DI PROCESSO
1923     init_descrittori_processi();
1924     // )
1925
1926     // ( CREAZIONE DEL PROCESSO MAIN_SISTEMA
1927     mid = crea_main_sistema();
1928     if (mid == 0xFFFFFFFF)
1929         goto error;
1929     flog(LOG_INFO, "Creato_il_processo_main_sistema_(id_=_%d)
1930         ↳ ", mid);
1931     // )
1932
1933     // ( CREAZIONE DEL PROCESSO DUMMY
1934     dummy_proc = crea_dummy();
1935     if (dummy_proc == 0xFFFFFFFF)
1936         goto error;
1936     flog(LOG_INFO, "Creato_il_processo_dummy_(id_=_%d)",
1937         ↳ dummy_proc);
1938     // )

```



```

1939 // (* SELEZIONIAMO MAIN_SISTEMA
1940 schedulatore();
1941 // *)
1942 // ( ESEGUE CALL CARICA_STATO; IRETQ (VEDI "SISTEMA.S")
1943 salta_a_main();
1944 // )
1945
1946 error:
1947     c_panic("Errore_di_inizializzazione");
1948 }

```

5.2 sistema.s

Listing 5.6: Salva e carica stato

```

38 // OFFSET DEI VARI REGISTRI ALL'INTERNO DI DES_PROC
39 .set punt_nucleo, 0
40 .set CR3, 8
41 .set RAX, CR3+8
42 .set RCX, CR3+16
43 .set RDX, CR3+24
44 .set RBX, CR3+32
45 .set RSP, CR3+40
46 .set RBP, CR3+48
47 .set RSI, CR3+56
48 .set RDI, CR3+64
49 .set R8, CR3+72
50 .set R9, CR3+80
51 .set R10, CR3+88
52 .set R11, CR3+96
53 .set R12, CR3+104
54 .set R13, CR3+112
55 .set R14, CR3+120
56 .set R15, CR3+128
57
58
59 //////////////////////////////////////
60 //          CAMBIO CONTESTO          //
61 //////////////////////////////////////
62
63 // COPIA LO STATO DEI REGISTRI GENERALI NEL DES_PROC DEL

```

```

64 // PROCESSO PUNTATO DA ESECUZIONE.
65 // NESSUN REGISTRO VIENE SPORCATO.
66 salva_stato:
67     // SALVIAMO LO STATO DI UN PAIO DI REGISTRI
68     // IN MODO DA POTERLI TEMPORANEAMENTE RIUTILIZZARE
69     // IN PARTICOLARE, USEREMO %RAX COME REGISTRO DI LAVORO
70     // E %RBX COME PUNTATORE AL DES_PROC.
71     .cfi_startproc
72     .cfi_def_cfa_offset 8
73     pushq %rbx
74     .cfi_adjust_cfa_offset 8
75     .cfi_offset rbx, -16
76     pushq %rax
77     .cfi_adjust_cfa_offset 8
78     .cfi_offset rax, -24
79
80     // RICAVIAMO IL PUNTATORE AL DES_PROC
81     movq esecuzione, %rax
82     movq $0, %rbx
83     movw (%rax), %bx    // CAMPO ID DEL PROC_ELEM
84     leaq descrittori_processi, %rax
85     movq (%rax, %rbx, 8), %rbx
86
87     // COPIAMO PER PRIMO IL VECCHIO VALORE DI %RAX
88     movq (%rsp), %rax
89     movq %rax, RAX(%rbx)
90     // USIAMO %RAX COME APPOGGIO PER COPIARE IL VECCHIO %RBX
91     movq 8(%rsp), %rax
92     movq %rax, RBX(%rbx)
93     // COPIAMO GLI ALTRI REGISTRI
94     movq %rcx, RCX(%rbx)
95     movq %rdx, RDX(%rbx)
96     // SALVIAMO IL VALORE CHE %RSP AVEVA PRIMA DELLA CHIAMATA
97     // A SALVA STATO (VALORE CORRENTE MENO GLI 8 BYTE CHE
98     // CONTENGONO L'INDIRIZZO DI RITORNO E I 16 BYTE DOVUTI
99     // ALLE DUE PUSH CHE ABBIAMO FATTO ALL'INIZIO)
100    movq %rsp, %rax
101    addq $24, %rax
102    movq %rax, RSP(%rbx)
103    movq %rbp, RBP(%rbx)
104    movq %rsi, RSI(%rbx)

```

```

105     movq %rdi, RDI(%rbx)
106     movq %r8, R8 (%rbx)
107     movq %r9, R9 (%rbx)
108     movq %r10, R10(%rbx)
109     movq %r11, R11(%rbx)
110     movq %r12, R12(%rbx)
111     movq %r13, R13(%rbx)
112     movq %r14, R14(%rbx)
113     movq %r15, R15(%rbx)
114
115     popq %rax
116     .cfi_adjust_cfa_offset -8
117     .cfi_restore rax
118     popq %rbx
119     .cfi_adjust_cfa_offset -8
120     .cfi_restore rbx
121
122     retq
123     .cfi_endproc
124
125
126 // CARICA NEI REGISTRI DEL PROCESSORE LO STATO CONTENUTO NEL
    ↪ DES_PROC DEL
127 // PROCESSO PUNTATO DA ESECUZIONE.
128 // QUESTA FUNZIONE SPORCA TUTTI I REGISTRI.
129 carica_stato:
130     .cfi_startproc
131     .cfi_def_cfa_offset 8
132
133     // CARICHIAMO NEL TSS IL PUNTATORE ALLA BASE DELLA PILA
    ↪ SISTEMA
134     // DEL PROCESSO IN ESECUZIONE
135     movq esecuzione, %rcx
136     movq $0, %rax
137     movw (%rcx), %ax          // ID PROCESSO
138     leaq descrittori_processi, %rbx
139     movq (%rbx, %rax, 8), %rbx // DES_PROC
140     movq punt_nucleo(%rbx), %rcx
141     leaq tss, %rax
142     movq %rcx, 4(%rax)
143

```

```

144     popq %rcx    //IND DI RITORNO, VA MESSO NELLA NUOVA PILA
145     .cfi_adjust_cfa_offset -8
146     .cfi_register rip, rcx
147
148     // NUOVO VALORE PER CR3
149     movq CR3(%rbx), %r10
150     movq %cr3, %rax
151     cmpq %rax, %r10
152     je 1f        // EVITIAMO DI INVALIDARE IL TLB
153                 // SE CR3 NON CAMBIA
154     movq %r10, %rax
155     movq %rax, %cr3    // IL TLB VIENE INVALIDATO
156 1:
157
158     // ANCHE SE ABBIAMO CAMBIATO CR3 SIAMO SICURI CHE
159     // L'ESECUZIONE PROSEGUE DA QUI, PERCHE' CI TROVIAMO
160     ↪ DENTRO
161     // LA FINESTRA FM CHE E' COMUNE A TUTTI I PROCESSI
162     movq RSP(%rbx), %rsp //CAMBIAMO PILA
163     pushq %rcx    //RIMETTIAMO L'INDIRIZZO DI RITORNO
164     .cfi_adjust_cfa_offset 8
165     .cfi_offset rip, -8
166
167     // SE IL PROCESSO PRECEDENTE ERA TERMINATO O ABORTITO LA
168     ↪ SUA PILA SISTEMA
169     // NON ERA STATA DISTRUTTA, IN MODO DA PERMETTERE A NOI
170     ↪ DI CONTINUARE
171     // AD USARLA. ORA CHE ABBIAMO CAMBIATO PILA POSSIAMO
172     ↪ DISFARCI DELLA
173     // PRECEDENTE.
174     cmpq $0, ultimo_terminato
175     je 1f
176     call distruggi_pila_precedente
177 1:
178
179     movq RCX(%rbx), %rcx
180     movq RDI(%rbx), %rdi
181     movq RSI(%rbx), %rsi
182     movq RBP(%rbx), %rbp
183     movq RDX(%rbx), %rdx
184     movq RAX(%rbx), %rax

```

```

181     movq R8(%rbx), %r8
182     movq R9(%rbx), %r9
183     movq R10(%rbx), %r10
184     movq R11(%rbx), %r11
185     movq R12(%rbx), %r12
186     movq R13(%rbx), %r13
187     movq R14(%rbx), %r14
188     movq R15(%rbx), %r15
189     movq RBX(%rbx), %rbx
190
191     retq
192     .cfi_endproc

```

Listing 5.7: Istruzioni per caricare nella GDT le primitive trampolino

```

333     carica_gate TIPO_TRAMP_IN    a_trampoline_in
           ↪ LIV_SISTEMA
334     carica_gate TIPO_TRAMP_OUT    a_trampoline_out
           ↪ LIV_SISTEMA

```

Listing 5.8: Inizializza GDT e TSS

```

364     .set p_dpl_type, 0b10001001 //P=1,DPL=00,TYPE=1001=TSS READY
365     .set pres_bit, 0b10000000
366
367     .global init_gdt
368     init_gdt:
369         lgdt gdt_pointer
370
371         // INIZIALIZZA DESCRITTORE DI TSS
372         leaq des_tss, %rdx
373         movw $DIM_DESP, (%rdx) // [15:0] = LIMIT [15:0]
374         decw (%rdx)
375         leaq tss, %rax
376         movw %ax, 2(%rdx) // [31:16] = BASE [15:0]
377         shrq $16, %rax
378         movb %al, 4(%rdx) // [39:32] = BASE [24:16]
379         movb $p_dpl_type, 5(%rdx) // [47:40] = P_DPL_TYPE
380         movb $0, 6(%rdx) // [55:48] = 0
381         movb %ah, 7(%rdx) // [63:56] = BASE [31:24]
382         shrq $16, %rax
383         movl %eax, 8(%rdx) // [95:64] = BASE [63:32]
384         movl $0, 12(%rdx) // [127:96] = 0

```

```

385
386 // CARICHIAMO TR CON L'OFFSET DELL'UNICO DESCRITTORE DI
    ↪ TSS
387 // (IN MODO CHE IL MECCANISMO DELLE INTERRUZIONI USI LA
388 // PILA SISTEMA DEL PROCESSO IN ESECUZIONE)
389 leaq des_tss, %rax
390 andb $0b11111101, 5(%rax) // RESET DEL BIT BUSY
391 // (RICHIESTO DAL PROCESSORE
392 // PER COMPATIBILITA' CON IL MODO
393 // A 32 BIT)
394 movq $(des_tss - gdt), %rax
395
396 ltr %ax
397
398 retq

```

Listing 5.9: Funzioni trampolino

```

519 //////////////////////////////////////
520 // SEZIONE TRAMPOLINE_TEXT: CODICE KERNEL //
521 //////////////////////////////////////
522 .section .trampoline_text, "ax"
523 .balign 4096
524 .global start_trampoline_text
525 start_trampoline_text:
526
527 // PARTE IN COMUNE TRA TRAMPOLINE_IN (CHIAMATA DAL MODULO
    ↪ SISTEMA)
528 // E A_TRAMPOLINE_IN (CHIAMATA DAL MODULO I/O)
529 trampoline_in_2:
530 .cfi_startproc
531 .cfi_def_cfa_offset 8
532
533 // MODIFICA LA TAB_ENTRY DELLA TABELLA DI LIVELLO 4
534 // RELATIVA ALLA FINESTRA DI MEMORIA FISICA ,
535 // SOSTITUENDO IL SOTTOALBERO DI TRADUZIONE SHADOW
536 // (USATO PER PROTEGGERE IL SISTEMA DA MELTDOWN)
537 // CON QUELLO KERNEL, COSI' DA POTER ACCEDERE ALLE
    ↪ FUNZIONALITA'
538 // DEL NUCLEO
539
540 movq des_finestra_kernel, %rbx

```

```

541     movq %cr3, %rax
542     movq %rbx, (%rax)
543
544     retq
545     .cfi_endproc
546
547 // PARTE IN COMUNE TRA TRAMPOLINE_OUT (CHIAMATA DAL MODULO
    ↪ SISTEMA)
548 // E A_TRAMPOLINE_OUT (CHIAMATA DAL MODULO I/O)
549 trampoline_out_2:
550     .cfi_startproc
551     .cfi_def_cfa_offset 8
552
553     // MODIFICA LA TAB_ENTRY DELLA TABELLA DI LIVELLO 4
554     // RELATIVA ALLA FINESTRA DI MEMORIA FISICA ,
555     // SOSTITUENDO IL SOTTOALBERO DI TRADUZIONE KERNEL
556     // CON QUELLO SHADOW, PER PROTEGGERE IL SISTEMA DA
    ↪ MELTDOWN.
557     // IN QUESTO MODO, IL PROGRAMMA UTENTE NON POTRA'
    ↪ ACCEDERE
558     // IN MANIERA SPECULATIVA ALLE ZONE DI MEMORIA DEL KERNEL
559
560     movq des_finestra_shadow, %rbx
561     movq %cr3, %rax
562     movq %rbx, (%rax)
563
564     // INVALIDIAMO IL TLB PER TOGLIERE LE TRADUZIONI DELLA
    ↪ FINESTRA DI MEMORIA FISICA
565     call     invalida_TLB
566
567     retq
568     .cfi_endproc
569
570
571 trampoline_in:
572     .cfi_startproc
573     .cfi_def_cfa_offset 8
574
575     pushq %rbx
576     .cfi_adjust_cfa_offset 8
577     .cfi_offset  rbx, -16

```

```

578     pushq %rax
579     .cfi_adjust_cfa_offset 8
580     .cfi_offset rax, -24
581
582     // CONTROLLIAMO SE STIAMO ENTRANDO DENTRO IL KERNEL DALLO
583     ↪ SPAZIO SHADOW,
584     // CONTROLLANDO A QUALE SOTTOSPAZIO APPARTIENE IL %RIP
585     ↪ SALVATO DALLA
586     // CHIAMATA DELL'INTERRUZIONE
587     movq 24(%rsp), %rbx    // %RIP SALVATO NELLA PILA DA INT
588     movabs $(1 << 47), %rax
589     addq %rbx, %rax
590     jnc 1f                // SALTA SE NON DOBBIAMO ENTRARE NEL
591     ↪ TRAMPOLINO
592
593     call trampoline_in_2
594
595 1: popq %rax
596     .cfi_adjust_cfa_offset -8
597     .cfi_restore rax
598     popq %rbx
599     .cfi_adjust_cfa_offset -8
600     .cfi_restore rbx
601     retq
602     .cfi_endproc
603
604 // EFFETTUA LA FASE DI TRAMPOLINO IN USCITA DAL KERNEL, DOPO
605 // AVER CONTROLLATO CHE IL PROGRAMMA TORNERA' A LIVELLO
606 ↪ UTENTE
607 trampoline_out:
608     .cfi_startproc
609     .cfi_def_cfa_offset 8
610
611     pushq %rbx
612     .cfi_adjust_cfa_offset 8
613     .cfi_offset rbx, -16
614     pushq %rax
615     .cfi_adjust_cfa_offset 8
616     .cfi_offset rax, -24
617
618     // CONTROLLIAMO SE STIAMO PASSANDO DA SISTEMA A UTENTE

```



```

615     movq 24(%rsp), %rbx      // %RIP SALVATO
616     movabs $(1 << 47), %rax
617     addq %rbx, %rax
618     jnc 1f                  // SALTA SE NON DOBBIAMO ENTRARE NEL
        ↪ TRAMPOLINO
619
620     // STIAMO USCENDO DAL KERNEL
621     // E' NECESSARIO IL TRAMPOLINO IN USCITA
622     call    trampoline_out_2
623
624 1:   popq %rax
625     .cfi_adjust_cfa_offset -8
626     .cfi_restore  rax
627     popq %rbx
628     .cfi_adjust_cfa_offset -8
629     .cfi_restore  rbx
630     retq
631     .cfi_endproc
632
633 // ROUTINE INT $TIPO_TRAMP_IN, ACCESSIBILE SOLO DAL MODULO I/
        ↪ O
634 a_trampoline_in:
635     .cfi_startproc
636     .cfi_def_cfa_offset 40
637     .cfi_offset rip, -40
638     .cfi_offset rsp, -16
639     pushq %rbx
640     .cfi_adjust_cfa_offset 8
641     .cfi_offset  rbx, -16
642     pushq %rax
643     .cfi_adjust_cfa_offset 8
644     .cfi_offset  rax, -24
645
646     call trampoline_in_2
647
648     popq %rax
649     .cfi_adjust_cfa_offset -8
650     .cfi_restore  rax
651     popq %rbx
652     .cfi_adjust_cfa_offset -8
653     .cfi_restore  rbx

```

```

654     iretq
655     .cfi_endproc
656
657 // ROUTINE INT $TIPO_TRAMP_OUT, ACCESSIBILE SOLO DAL MODULO I
    ↪ /O
658 a_trampoline_out:
659     .cfi_startproc
660     .cfi_def_cfa_offset 40
661     .cfi_offset rip, -40
662     .cfi_offset rsp, -16
663     pushq %rbx
664     .cfi_adjust_cfa_offset 8
665     .cfi_offset rbx, -16
666     pushq %rax
667     .cfi_adjust_cfa_offset 8
668     .cfi_offset rax, -24
669
670     call trampoline_out_2
671
672     popq %rax
673     .cfi_adjust_cfa_offset -8
674     .cfi_restore rax
675     popq %rbx
676     .cfi_adjust_cfa_offset -8
677     .cfi_restore rbx
678     iretq
679     .cfi_endproc

```

Listing 5.10: Esempio della parte Assembly di una routine di sistema
con le chiamate alle funzioni trampolino

```

685 a_activate_p: // ROUTINE INT $TIPO_A
686     .cfi_startproc
687     .cfi_def_cfa_offset 40
688     .cfi_offset rip, -40
689     .cfi_offset rsp, -16
690     call trampoline_in
691     call salva_stato
692     cavallo_di_troia %rdi
693     call c_activate_p
694     call carica_stato
695     call trampoline_out

```

```

696     iretq
697     .cfi_endproc

```

Listing 5.11: Sezioni DATA e TSS della memoria shadow/trampolino

```

1666 .global end_trampoline_text
1667 end_trampoline_text:
1668
1669 //////////////////////////////////////
1670 //      SEZIONE TRAMPOLINE_DATA: TABELLE E DATI      //
1671 //////////////////////////////////////
1672 .section .trampoline_data, "aw"
1673 .balign 4096
1674 .global start_trampoline_data
1675 start_trampoline_data:
1676
1677 // PUNTATORI ALLE TABELLE GDT E IDT
1678 // NEL FORMATO RICHIESTO DALLE ISTRUZIONI LGDT E LIDT
1679 gdt_pointer:
1680     .word end_gdt-gdt      // LIMITE DELLA GDT
1681     .quad gdt              // BASE DELLA GDT
1682 idt_pointer:
1683     .word 0xFFF           // LIMITE DELLA IDT (256 ENTRATE)
1684     .quad idt             // BASE DELLA IDT
1685 triple_fault_idt:
1686     .word 0
1687     .quad 0
1688 param_err:
1689     .asciz "indirizzo non valido: %p"
1690
1691 .balign 8
1692 .global gdt
1693 gdt:
1694     .quad 0              //SEGMENTO NULLO
1695 code_sys_seg:
1696     .word 0b0            //LIMIT[15:0]    NOT USED
1697     .word 0b0            //BASE[15:0]     NOT USED
1698     .byte 0b0            //BASE[23:16]    NOT USED
1699     .byte 0b10011010     //P|DPL|1|1|C|R|A|  DPL=00=SISTEMA
1700     .byte 0b00100000     //G|D|L|-----|  L=1 LONG MODE
1701     .byte 0b0            //BASE[31:24]    NOT USED
1702 code_usr_seg:

```

```

1703     .word 0b0          //LIMIT[15:0]    NOT USED
1704     .word 0b0          //BASE[15:0]      NOT USED
1705     .byte 0b0          //BASE[23:16]     NOT USED
1706     .byte 0b11111010   //P|DPL|1|1|C|R|A|   DPL=11=UTENTE
1707     .byte 0b00100000   //G|D|L|-|-----|   L=1 LONG MODE
1708     .byte 0b0          //BASE[31:24]     NOT USED
1709 data_usr_seg:
1710     .word 0b0          //LIMIT[15:0]    NOT USED
1711     .word 0b0          //BASE[15:0]      NOT USED
1712     .byte 0b0          //BASE[23:16]     NOT USED
1713     .byte 0b11110010   //P|DPL|1|0|E|W|A|   DPL=11=UTENTE
1714     .byte 0b00000000   //G|D|-|-|-----|
1715     .byte 0b0          //BASE[31:24]     NOT USED
1716 des_tss:
1717     .space 16,0 //SEGMENTO TSS , RIEMPITO A RUNTIME
1718 end_gdt:
1719
1720 .global tss
1721 tss:
1722     .fill DIM_DESP, 0
1723
1724 .global des_finestra_kernel
1725 des_finestra_kernel:
1726     .quad 0
1727
1728 .global end_trampoline_data
1729 end_trampoline_data:
1730
1731
1732 .section .trampoline_bss, "aw", @nobits
1733 .balign 4096
1734 .global start_trampoline_bss
1735 start_trampoline_bss:
1736
1737 // .BALIGN 16
1738 idt:
1739     // SPAZIO PER 256 GATE
1740     // VERRA' RIEMPITA A TEMPO DI ESECUZIONE
1741     .space 16 * 256, 0
1742 end_idt:
1743

```

```

1744 exc_error:
1745     .space 8, 0
1746
1747 .global end_trampoline_bss
1748 end_trampoline_bss:

```

5.3 io.s

Listing 5.12: Esempio della parte Assembly di una primitiva di I/O
con le chiamate alle primitive trampolino

```

378 a_readse_n:
379     .cfi_startproc
380     .cfi_def_cfa_offset 40
381     .cfi_offset rip, -40
382     .cfi_offset rsp, -16
383     cavallo_di_troia %rsi
384     cavallo_di_troia2 %rsi %rdx
385     cavallo_di_troia %rcx
386     int $TIPO_TRAMP_IN
387     call c_readse_n
388     int $TIPO_TRAMP_OUT
389     iretq
390     .cfi_endproc

```

5.4 costanti.h

Listing 5.13: Costanti per gli id delle primitive trampoline

```

76 #define TIPO_TRAMP_IN      0x59    // TRAMPOLINE IN  INGRESSO
77 #define TIPO_TRAMP_OUT    0x5a    // TRAMPOLINE IN  USCITA

```


Listings

| | | |
|------|--|----|
| 2.1 | Meltdown frame | 10 |
| 5.1 | Descrittori di processo e di TSS | 19 |
| 5.2 | Creazione della finestra di memoria shadow | 21 |
| 5.3 | Creazione di un processo | 26 |
| 5.4 | Distruzione di un processo | 29 |
| 5.5 | Inizializzazione del sistema | 30 |
| 5.6 | Salva e carica stato | 33 |
| 5.7 | Istruzioni per caricare nella GDT le primitive trampolino | 37 |
| 5.8 | Inizializza GDT e TSS | 37 |
| 5.9 | Funzioni trampolino | 38 |
| 5.10 | Esempio della parte Assembly di una routine di sistema con le chiamate alle funzioni trampolino | 42 |
| 5.11 | Sezioni DATA e TSS della memoria shadow/trampolino | 43 |
| 5.12 | Esempio della parte Assembly di una primitiva di I/O con le chiamate alle primitive trampolino | 45 |
| 5.13 | Costanti per gli id delle primitive trampoline | 45 |

Bibliografia

- [1] Graziano Frosini e Giuseppe Lettieri. *Architettura dei calcolatori Vol. II: Struttura hardware del processore PC, del Bus, della memoria, delle interfacce e gestione dell'I/O*. A cura di Pisa University Press. 2013.
- [2] Graziano Frosini e Giuseppe Lettieri. *Architettura dei calcolatori Vol. III: Aspetti architetturali avanzati e nucleo di sistema operativo*. A cura di Pisa University Press. 2013.
- [3] Daniel Gruss et al. «KASLR is Dead: Long Live KASLR». In: *International Symposium on Engineering Secure Software and Systems*. Springer International Publishing, 2017, pp. 161–176.
- [4] Daniel Gruss et al. «Prefetch Side-Channel Attacks: Bypassing SMAP and Kernel ASLR». In: *Conference on Computer and Communications Security*. 2016, pp. 368–379.
- [5] Ralf Hund, Carsten Willems e Thorsten Holz. «Practical Timing Side Channel Attacks against Kernel Space ASLR». In: *Security and Privacy*. 2013, pp. 191–205.
- [6] Yeongjin Jang, Sangho Lee e Taesoo Kim. «Breaking Kernel Address Space Layout Randomization with Intel TSX». In: *Conference on Computer and Communications Security*. 2016, pp. 380–392.
- [7] Giuseppe Lettieri. *Implementazione della memoria virtuale*. 13 Mag. 2019. URL: <https://calcolatori.iet.unipi.it/resources/paginazione-nel-nucleo.pdf>.
- [8] Giuseppe Lettieri. *Introduzione al sistema multiprogrammato*. 24 Apr. 2020. URL: <https://calcolatori.iet.unipi.it/resources/nucleo.pdf>.
- [9] Giuseppe Lettieri. *Istruzioni all'uso del nucleo*. URL: https://calcolatori.iet.unipi.it/istruzioni_nucleo.php.
- [10] Giuseppe Lettieri. *Memoria Cache*. 16 Mar. 2017. URL: <https://calcolatori.iet.unipi.it/resources/cache.pdf>.

- [11] Giuseppe Lettieri. *Paginazione*. 3 Mag. 2019. URL: <https://calcolatori.iet.unipi.it/resources/paginazione.pdf>.
- [12] Giuseppe Lettieri. *Paginazione su domanda*. 14 Mag. 2018. URL: <https://calcolatori.iet.unipi.it/resources/paginazione-su-domanda.pdf>.
- [13] Giuseppe Lettieri. *Paginazione: complementi*. 4 Apr. 2017. URL: <https://calcolatori.iet.unipi.it/resources/tlb.pdf>.
- [14] Giuseppe Lettieri. *Protezione*. 11 Apr. 2019. URL: <https://calcolatori.iet.unipi.it/resources/protezione.pdf>.
- [15] Moritz Lipp et al. «Meltdown: Reading Kernel Memory from User Space». In: *USENIX Security Symposium*. 2018.
- [16] Yuval Yarom e Katrina Falkner. «FLUSH+RELOAD: A High Resolution, Low Noise, L3 Cache Side-Channel Attack». In: *USENIX Security Symposium*. 2014, pp. 719–732.